

# Middlesex University Research Repository

An open access repository of

Middlesex University research

<http://eprints.mdx.ac.uk>

Blanchette, Jasmin Christian and Popescu, Andrei (2013) Mechanizing the metatheory of sledgehammer. *Frontiers of Combining Systems: 9th International Symposium, FroCoS 2013, Nancy, France, September 18-20, 2013. Proceedings.* In: *9th International Symposium on Frontiers of Combining Systems (FroCoS 2013), 18-20 Sept 2013, Nancy, France.* ISBN 9783642408847. ISSN 0302-9743 [Conference or Workshop Item]

Final accepted version (with author's formatting)

This version is available at: <https://eprints.mdx.ac.uk/15370/>

## Copyright:

Middlesex University Research Repository makes the University's research available electronically.

Copyright and moral rights to this work are retained by the author and/or other copyright owners unless otherwise stated. The work is supplied on the understanding that any use for commercial gain is strictly forbidden. A copy may be downloaded for personal, non-commercial, research or study without prior permission and without charge.

Works, including theses and research projects, may not be reproduced in any format or medium, or extensive quotations taken from them, or their content changed in any way, without first obtaining permission in writing from the copyright holder(s). They may not be sold or exploited commercially in any format or medium without the prior written permission of the copyright holder(s).

Full bibliographic details must be given when referring to, or quoting from full items including the author's name, the title of the work, publication details where relevant (place, publisher, date), pagination, and for theses or dissertations the awarding institution, the degree type awarded, and the date of the award.

If you believe that any material held in the repository infringes copyright law, please contact the Repository Team at Middlesex University via the following email address:

[eprints@mdx.ac.uk](mailto:eprints@mdx.ac.uk)

The item will be removed from the repository while any claim is being investigated.

See also repository copyright: re-use policy: <http://eprints.mdx.ac.uk/policies.html#copy>

# Mechanizing the Metatheory of Sledgehammer

Jasmin Christian Blanchette and Andrei Popescu

Fakultät für Informatik, Technische Universität München, Germany

**Abstract.** This paper presents an Isabelle/HOL formalization of recent research in automated reasoning: efficient encodings of sorts in unsorted first-order logic, as implemented in Isabelle’s Sledgehammer proof tool. The formalization provides the general-purpose machinery to reason about formulas and models, emulating the theory of institutions. Quantifiers are represented using a nominal-like approach designed for interpreting syntax in semantic domains.

## 1 Introduction

Despite steady progress in the usability of proof assistants, paper proofs reign supreme in the automated reasoning community. Myreen and Davis’s verification of an ACL2-like prover in HOL4 [16] and Harrison’s partial self-verification of HOL Light [12] are exceptions rather than the rule. Important metamathematical results have been formalized (e.g., Shankar’s Gödel proof [26]), but new research is still carried out almost exclusively on paper, with all the risks this entails.

This paper presents a formalization in Isabelle/HOL [17] of the proofs for translations from many-sorted to unsorted first-order logic (FOL). Claessen et al. [9] designed lightweight encodings that eliminate much of the clutter associated with traditional schemes. Blanchette et al. [3, 4] introduced even lighter encodings in a sequel. Central to these new encodings is the notion of monotonicity. Informally, a sort is monotonic if its domain can be extended with new elements without compromising satisfiability. Nonmonotonic sorts can be made monotonic by introducing protector functions or predicates, and monotonic sorts can be merged into a single sort.

Sorts are trivially monotonic in FOL without equality. The addition of interpreted equality makes it possible to encode upper cardinality bounds on the models, breaking monotonicity. Like other interesting semantic properties, monotonicity is undecidable but can often be inferred in practice. Monotonicity has many applications in theorem provers and model finders [5, 9]. It is also roughly equivalent to smoothness, a criterion that arises when combining decision procedures in SMT solvers [28].

The Sledgehammer [18] proof tool for Isabelle/HOL relies on the monotonicity-based encodings to apply state-of-the-art unsorted provers to sorted problems. The tool translates interactive proof goals along with relevant lemmas and invokes the external automatic theorem provers to find proofs, which are reconstructed through Isabelle’s inference kernel. Early versions of Sledgehammer relied on unsound sort encodings; as a result, they would often find spurious, unreconstructable proofs, which annoyed users and could conceal sound proofs. Whereas Sledgehammer reconstructs the external proofs, tools such as Monotonox [9] and the fully-automatic competition version of Isabelle [27] do not perform such checks; soundness is crucial for them.

The mechanization of the sort encodings fully covers the correctness proofs from Claessen et al. [9] and the monomorphic half of its sequel [3, 4], as well as a theorem by Bouillaguet et al. [8]. This formalization work arose from a desire to provide more solid assurance to this recent research. Even if the intuition is clear, a paper proof offers many opportunities for flaws, especially because of the variety of encodings.

The mechanization effort partly coincided with the development of the informal proofs [4]. The two proofs largely follow the same conventions, with one major difference: The core of the formal proof (Sections 3 to 5) assumes quantifier-free clausal normal form (CNF) rather than negation normal form (NNF). This reduces the exposure to name binders, which are notoriously difficult to reason about. The results are lifted to NNF using a clausification theorem (Sections 6 to 8). This organization is reminiscent of the architecture of automatic reasoners that combine a clausifier and a CNF core.

Isabelle’s higher-order logic (HOL) might not be as expressive as set or type theory, but it can cope with the statements and proofs of classical metatheorems (as shown by Harrison and others [2, 11, 25]) and practical results. The proof assistant offers many conveniences; two features have been particularly useful:

- *Locales* [1, 14] parameterize theories over constants and assumptions, with the usual benefits associated with modularity. Locales are particularly suited to expressing logic translations abstractly as in the theory of institutions [10].
- A framework for *syntax with bindings* [22–24] eases reasoning about quantified formulas. It lies at the intersection of first-order nominal approaches [20] and higher-order abstract syntax [19]. The framework is designed specifically for interpreting syntax in semantic domains.

Locales have been part of Isabelle for many years and are widely used. The syntax with bindings is a newer addition; the current application is among the first case studies that feature it. The formal proofs are available online [6].

Although sort encodings are the focus of this paper, our infrastructure is designed to be reusable for other applications of many-sorted FOL. Many important metatheories are awaiting formalization, such as the completeness of paramodulation and tableaux.

## 2 An Isabelle View of Logic Translations

The formalization covers a variety of translations, including not only the sort encodings but also clausification. The guiding principles, described below, originate from the theory of institutions; their Isabelle materialization relies on locales.

**Institutions.** A logic  $\mathcal{L}$  provides a category of signatures  $Sig$  and, for each signature  $\Sigma \in Sig$ , a set of sentences  $Sen(\Sigma)$ , a class of structures (interpretations)  $Str(\Sigma)$ , and a satisfaction relation  $\models_\Sigma$  between structures and sentences. A signature morphism  $k : \Sigma \rightarrow \Sigma'$  is equipped with a forward sentence translation  $k : Sen(\Sigma) \rightarrow Sen(\Sigma')$  and a backward structure translation  $\downarrow_k : Str(\Sigma') \rightarrow Str(\Sigma)$ . An *institution* is a logic whose signature morphisms enjoy the property that “truth is invariant under change of notation”:  $\mathcal{M}' \models_{\Sigma'} k \varphi \iff \mathcal{M}' \downarrow_k \models_\Sigma \varphi$  for all  $k : \Sigma \rightarrow \Sigma'$ ,  $\mathcal{M}' \in Str(\Sigma')$ , and  $\varphi \in Sen(\Sigma)$ .

A translation of  $\mathcal{L}$ -problems (sets of sentences) into  $\mathcal{L}'$ -problems consists of a function  $\$$  between  $\mathcal{L}$ 's and  $\mathcal{L}'$ 's signature classes and, for each  $\Sigma \in \text{dom}(\$)$  and  $\Sigma$ -problem  $\Phi$ , a sentence translation  $\text{enc}_\Phi : \text{Sen}(\Sigma) \rightarrow \text{Sen}(\Sigma^\$)$  and a set of axioms  $\mathcal{A}x_\Phi \subseteq \text{Sen}(\Sigma^\$)$ . The translation of  $\Phi$  is defined as  $\text{enc } \Phi = \{\text{enc}_\Phi \varphi \mid \varphi \in \Phi\} \cup \mathcal{A}x_\Phi$ . Thus,  $\mathcal{L}$ -problems are mapped to  $\mathcal{L}'$ -problems by joining an elementwise translation and additional axioms. Given a class  $C$  of  $\mathcal{L}$ -problems, the translation is *sound* w.r.t.  $C$  if satisfiability of  $\Phi$  implies satisfiability of  $\text{enc } \Phi$  for all  $\Phi \in C$ , and *complete* if the converse holds.

The institution literature focuses on “uniform” encodings. For these, the sentence translation depends only on  $\Phi$ 's signature  $\Sigma$ , and there exists a backward translation  $\text{dec} : \text{Str}(\Sigma^\$) \rightarrow \text{Str}(\Sigma)$  for which an inter-institution version of the institutional condition holds:  $\mathcal{M}' \models_{\Sigma^\$} \text{enc}_\Sigma \varphi \iff \text{dec } \mathcal{M}' \models_\Sigma \varphi$ . This condition implies completeness.

The source logic  $\mathcal{L}$  for all the translations considered in this paper is many-sorted FOL; the target logic  $\mathcal{L}'$  is either many-sorted or unsorted FOL. Sentences are either CNF clauses or NNF formulas. Most of the translations are nonuniform.

**Isabelle.** Isabelle/HOL is based on polymorphic HOL, which can be thought of as a fragment of Standard ML enriched with logical constructs and a proof system. Type variables are identified by a leading prime (e.g.,  $'a$ ). The type  $\sigma \rightarrow \tau$  is interpreted as the set of (total) functions from  $\sigma$  to  $\tau$ . Propositions are terms of type `bool`, and predicates are functions to `bool`. Function applications are written without parentheses (e.g.,  $f \ x \ y$ ) or in infix notation (e.g.,  $x + y$ ). Constants and variables can be functions.

The type  $'a$  list of finite lists over  $'a$  is generated freely from the empty list `[]` and the infix constructor `#` :  $'a \rightarrow 'a \text{ list} \rightarrow 'a \text{ list}$ . The notation  $[x_1, x_2, \dots, x_n]$  abbreviates  $x_1 \# (x_2 \# (\dots \# (x_n \# []) \dots))$ . The higher-order constant `map` :  $('a \rightarrow 'b) \rightarrow 'a \text{ list} \rightarrow 'b \text{ list}$  applies a unary function to each element in a list, and `set` :  $'a \text{ list} \rightarrow 'a \text{ set}$  returns the set of elements in a list. Sets are written using traditional mathematical notation. Type parameters of polymorphic types are sometimes omitted (e.g., `set` for  $'a \text{ set}$ ).

**Locales.** Isabelle locales are a structuring mechanism provided on top of basic HOL. They fix types, constants, and assumptions, as in the following schematic examples:

```
locale X = fixes 'a fixes c :  $\sigma_{'a}$  assumes  $P_{'a,c}$ 
locale Y = fixes 'b fixes d :  $\tau_{'b}$  assumes  $Q_{'b,d}$ 
```

The definition of locale `X` fixes a type  $'a$  and a constant `c` whose type  $\sigma_{'a}$  may depend on  $'a$ , and states an assumption  $P_{'a,c} : \text{bool}$  over  $'a$  and `c`. Lemmas proved within the locales can rely on them. In general, a single locale can introduce several types, constants, and assumptions. The definition of `X` also produces a polymorphic *locale predicate*  $X = (\lambda c. P_{'a,c})$ . Seen from outside the locale, the lemmas proved in locale `X` are polymorphic in type variable  $'a$ , universally quantified over variable  $c$ , and conditional on  $X \ c$ .

Locales support inheritance, union, and embedding. To embed `X` into `Y`, one needs to indicate how an arbitrary instance of `X` can be regarded as an instance of `Y`, by providing, in the context of `X`, definitions of the types and constants of `Y` together with proofs of `Y`'s assumptions. The command

```
sublocale X < Y where 'b =  $\nu$  and d =  $t$ 
```

emits the goal  $Q_{\nu,t}$ , where  $\nu$  and  $t$  may depend on types and constants from `X`. After the proof, all the lemmas proved in the `Y` become available in `X`, with  $\nu$  and  $t$  in place of  $'b$  and `d`. Homonymous constants `d` in `X` and `Y` are instantiated as  $d = d$  by default.

The sublocale relationship is sometimes abbreviated to  $X'_{a,c} < Y'_{u,t}$  or  $X < Y$ .

Locales provide a shallow realization of institutions in Isabelle. The institutional methodology serves as an inspiration and guidance to formulate results about specific logic translations in a consistent style. Given a logic  $\mathcal{L}$ , its signatures  $Sig$  are captured by a locale  $\mathcal{L}.Signature$ , which fixes Isabelle constants for the signature components (e.g., sorts and symbols) and defines a notion of sentence (e.g., clauses or formulas). A locale  $\mathcal{L}.Problem$  extends  $\mathcal{L}.Signature$  with a fixed set of sentences  $\Phi$ . Structures  $\mathcal{M}$  are represented by a locale  $\mathcal{L}.Structure$  that also defines a notion of satisfaction. Finally, satisfiable problems are represented by a locale  $\mathcal{L}.Model$  that joins  $\mathcal{L}.Problem$  and  $\mathcal{L}.Structure$  and further requires satisfaction between  $\Phi$  and  $\mathcal{M}$ .

In this setting, translations between logics  $\mathcal{L}$  and  $\mathcal{L}'$  and their properties are captured via locale embedding mechanisms in four steps.

**SIG:** Define  $\$$  as a sublocale relationship  $\mathcal{L}.Signature < \mathcal{L}'.Signature$  with suitable parameter instantiations reflecting the definition of  $\Sigma^\$$  in terms of  $\Sigma$ .

**TRANS:** Define  $enc_\Phi$  inside  $\mathcal{L}.Problem$  (where  $\Sigma$  and the  $\Sigma$ -problem  $\Phi$  are fixed).

**SOUND:** To prove soundness, define a  $\Sigma^\$$ -structure  $\mathcal{M}'$  inside  $\mathcal{L}.Model$  (where the signature  $\Sigma$ , the  $\Sigma$ -problem  $\Phi$ , and the structure  $\mathcal{M}$  such that  $\Phi \models_\Sigma \mathcal{M}$  are fixed) and show  $\mathcal{L}.Model_{\mathcal{M}} < \mathcal{L}'.Model_{\mathcal{M}'}$ .

**COMPLETE:** To prove completeness, define a locale  $Problem\_Model' = \mathcal{L}.Problem + \mathcal{L}'.Model$  that joins a  $\Sigma$ -problem  $\Phi$  and a  $\Sigma^\$$ -model  $\mathcal{M}'$  of  $enc_\Phi$ , define inside  $Problem\_Model'$  a  $\Sigma$ -structure  $\mathcal{M}$ , and show  $Problem\_Model'_{\mathcal{M}'} < \mathcal{L}.Model_{\mathcal{M}}$ .

### 3 Clausal First-Order Logic

The terms, atoms, and literals of (quantifier-free) CNF are represented in HOL by ML-style free datatypes, parameterized by types  $'f$  and  $'p$  of function and predicate symbols:

<code>datatype 'f tm =</code>	<code>datatype ('f,'p) atm =</code>	<code>datatype ('f,'p) lit =</code>
<code>Var var  </code>	<code>Pr 'p ('f tm list)  </code>	<code>Pos (('f,'p) atm)  </code>
<code>Fm 'f ('f tm list)</code>	<code>Eq ('f tm) ('f tm)</code>	<code>Neg (('f,'p) atm)</code>

The type `var` is countably infinite. An atom is either an applied predicate (e.g.,  $p(t)$ ) or equality (e.g.,  $t \approx u$ ). A clause is a list of literals (interpreted disjunctively), and a problem is a set of clauses (interpreted conjunctively). Formally,  $(f,p)$  clause =  $(f,p)$  lit list and  $(f,p)$  problem =  $(f,p)$  clause set. The CNF representation involves no name binders, unlike (quantified) NNF (Section 6).

Many-sorted signatures (for CNF and NNF) are captured by the following locale:

```
locale Signature =
  fixes 's and 'f and 'p
  fixes arity_F : 'f → 's list and res : 'f → 's and arity_P : 'p → 's list
  assumes countable UNIV'_s and countable UNIV'_f and countable UNIV'_p
```

The locale is parameterized by types for sorts ( $'s$ ), function symbols ( $'f$ ), and predicate symbols ( $'p$ ), all required to be countable (i.e. finite or countably infinite). The locale attaches to each symbol a sort `arity` (`arity_F` or `arity_P`) and, for functions, a result sort (`res`). The sort `arity` can be empty. Symbols cannot be overloaded. The polymorphic constant `UNIV'_a : 'a set` is predefined in Isabelle as the set of all values of type  $'a$ .

The Signature locale defines an underspecified function  $\text{sort} : \text{var} \rightarrow 's$  that arbitrarily assigns sorts to variables. Whereas the formalization consistently refers to FOL's sorts as types (in view of a possible extension to  $n$ -ary type constructors and polymorphism), in this paper they are more precisely called sorts. Wellsortedness and wellformedness of terms and the other syntactic categories are defined in the usual way. Wellformedness is a precondition to many operations, but such details are omitted here.

The Problem locale joins a signature  $\Sigma$  and a CNF  $\Sigma$ -problem  $\Phi$ . The Structure locale combines a signature, a universe  $'u$ , and a triple of functions  $(\text{int}_S, \text{int}_F, \text{int}_P)$  that interpret sorts, function symbols, and predicate symbols:

```

locale Problem = Signature_{s,f,p} arity_F res arity_P +
  fixes  $\Phi : ('f, 'p)$  problem

locale Structure = Signature_{s,f,p} arity_F res arity_P +
  fixes  $'u$ 
  fixes  $\text{int}_S : 's \rightarrow 'u \rightarrow \text{bool}$  and  $\text{int}_F : 'f \rightarrow 'u \text{ list} \rightarrow 'u$  and
   $\text{int}_P : 'p \rightarrow 'u \text{ list} \rightarrow \text{bool}$ 

```

A few wellformedness assumptions are made on the triple  $(\text{int}_S, \text{int}_F, \text{int}_P)$ , such as inhabitation of all sorts  $(\forall \sigma. \exists d. \text{int}_S \sigma d)$ . The Structure locale also defines the interpretation of terms and satisfaction of clauses. A related locale, Model, represents satisfiable CNF problems by combining a Problem and a Structure it satisfies.

## 4 Monotonicity and Its Inference

This section focuses on monotonicity in its own right; Section 5 discusses the associated sort encodings. To simplify the monotonicity arguments, both sections assume a fixed infinitely countable type  $\omega$  as the universe  $'u$  of structures, thus working implicitly with the instances  $\text{Structure}_\omega$  and  $\text{Model}_\omega$ . This limitation is lifted in Section 8 by appealing to the downward Löwenheim–Skolem theorem.

Claessen et al. [9, §2] define monotonicity on single sorts. Blanchette et al. [3, §3] generalized the notion to sets of sorts  $S$ , making it more useful. The sorts  $S$  are collectively *monotonic* in the problem  $\Phi$  if for all models  $\mathcal{M}$  of  $\Phi$ , there exists a model  $\mathcal{M}'$  such that for all sorts  $\sigma$ ,  $\mathcal{M}'$  interprets  $\sigma$  by an infinite domain if  $\sigma \in S$  and by a domain of the same cardinality as in  $\mathcal{M}$  otherwise.

In the formalization, the Mono\_Problem locale enriches Problem with a monotonicity assumption on all sorts, expressed using locale predicates:

$$(\exists \text{int}_S \text{int}_F \text{int}_P. \text{Model } \text{arity}_F \text{ res } \text{arity}_P \Phi \text{ int}_S \text{int}_F \text{int}_P) \longrightarrow \exists \text{int}_S \text{int}_F \text{int}_P. \text{Infinite\_Model } \text{arity}_F \text{ res } \text{arity}_P \Phi \text{ int}_S \text{int}_F \text{int}_P$$

The Infinite\_Model locale is itself an enrichment of Model with the assumption that for each sort  $\sigma$ , the expression  $\text{int}_S \sigma d$  is true for infinitely many elements  $d$ .

**First Criterion.** Claessen et al. designed two syntactic criteria to infer monotonicity. The first one is defined as a predicate  $\triangleright$  that checks the absence of naked variables of a given sort  $\sigma$  in a clause  $c$  or a problem  $\Phi$ :

$$\sigma \triangleright c \longleftrightarrow \forall x \in \text{nv } c. \text{sort } x \neq \sigma \quad \sigma \triangleright \Phi \longleftrightarrow \forall c \in \Phi. \sigma \triangleright c$$

A *naked variable* is a variable that occurs directly on either side of a positive equality, such as  $X$  in the literal  $X \approx f(Y)$ . Formally:

$$\begin{array}{lll} \text{nv}(\text{Var } x) = \{x\} & \text{nv}(\text{Eq } t_1 t_2) = \text{nv } t_1 \cup \text{nv } t_2 & \text{nv}(\text{Pos } a) = \text{nv } a \\ \text{nv}(\text{Fn } f ts) = \emptyset & \text{nv}(\text{Pr } p ts) = \emptyset & \text{nv}(\text{Neg } a) = \emptyset \end{array}$$

with  $\text{nv } c = \bigcup \text{set}(\text{map } \text{nv } c)$  for clauses. The criterion  $\triangleright$  soundly infers monotonicity. This is expressed as a sublocale inclusion  $\text{Problem\_Crit1} < \text{Mono\_Problem}$ , where  $\text{Problem\_Crit1}$  enriches  $\text{Problem}$  with the assumption  $\forall \sigma. \sigma \triangleright \Phi$ . The inclusion holds because a model of a problem whose sorts pass  $\triangleright$  can be extended into an infinite model by replicating elements. For each finite sort  $\sigma$ , the extended model contains infinitely many copies of some element pick  $\sigma$ , all interpreted as in the original model.

Blanchette et al. strengthened the criterion by injecting “infinity knowledge”: Any sort that is interpreted by an infinite domain in all models is monotonic, regardless of naked variables [3, §3]. This aspect is part of the formalization but omitted here.

**Second Criterion.** The improved criterion is parameterized by an assignment of a per-sort *extension policy*—which may be *true*, *false*, or *copy*—to each predicate symbol. In the model construction, the true-extended (resp. false-extended) predicates are interpreted as true (resp. false) for new domain elements of the given sort, whereas the copy-extended predicates are treated as in the simple criterion.

Implementations can enumerate the possible policy combinations (e.g., using a SAT solver). In the formalization, the policies are supplied along with the problem as a curried function *policy* that maps pairs  $\sigma, p$  to T, F, or C. A function *guard* associates each variable  $x$  in need of protection with its guarding literal. The criterion is defined as

$$\begin{array}{l} \sigma \triangleright c \longleftrightarrow \forall l x. l \in \text{set } c \wedge x \in \text{nv } l \wedge \text{sort } x = \sigma \longrightarrow \text{isGuard } x (\text{guard } c \ l \ x) \\ \sigma \triangleright \Phi \longleftrightarrow \forall c \in \Phi. \sigma \triangleright c \end{array}$$

where *isGuard* determines whether the given literal actually protects the variable  $x$ :

$$\begin{array}{l} \text{isGuard } x (\text{Pos } (\text{Eq } t_1 t_2)) \longleftrightarrow \text{False} \\ \text{isGuard } x (\text{Neg } (\text{Eq } t_1 t_2)) \longleftrightarrow \bigvee_{i=1}^2 t_i = \text{Var } x \wedge \exists f ts. t_{3-i} = \text{Fn } f ts \\ \text{isGuard } x (\text{Pos } (\text{Pr } p ts)) \longleftrightarrow x \in \bigcup \text{set}(\text{map } \text{nv } ts) \wedge \text{policy}(\text{sort } x) p = \text{T} \\ \text{isGuard } x (\text{Neg } (\text{Pr } p ts)) \longleftrightarrow x \in \bigcup \text{set}(\text{map } \text{nv } ts) \wedge \text{policy}(\text{sort } x) p = \text{F} \end{array}$$

The notion of naked variables is generalized to account for ill-polarized predicates:

$$\begin{array}{l} \text{nv}(\text{Pos } (\text{Pr } p ts)) = \{x \in \bigcup \text{set}(\text{map } \text{nv } ts) \mid \text{policy}(\text{sort } x) p = \text{F}\} \\ \text{nv}(\text{Neg } (\text{Pr } p ts)) = \{x \in \bigcup \text{set}(\text{map } \text{nv } ts) \mid \text{policy}(\text{sort } x) p = \text{T}\} \end{array}$$

**Theorem 1.** Let  $\Phi$  be a  $\Sigma$ -problem and  $\sigma$  be a  $\Sigma$ -sort.

- (1) If  $\sigma \triangleright \Phi$ , then  $\sigma \triangleright \Phi$  for a copy-extended policy.
- (2) Given some extension policies, if  $\sigma \triangleright \Phi$  for all  $\Sigma$ -sorts  $\sigma$ , then the set of all  $\Sigma$ -sorts is monotonic in  $\Phi$ .

This theorem is expressed in Isabelle as a pair of sublocale inclusions. The *where* clause below instantiates  $\text{Problem\_Policy\_Crit2}$ ’s *policy* parameter with  $\lambda \sigma p. \text{C}$  to enforce the copy policy for all sorts and predicate symbols:

$$\begin{array}{l} \text{sublocale } \text{Problem\_Crit1} < \text{Problem\_Policy\_Crit2} \text{ where } \text{policy} = (\lambda \sigma p. \text{C}) \\ \text{sublocale } \text{Problem\_Policy\_Crit2} < \text{Mono\_Problem} \end{array}$$

## 5 Sort Encodings

A naive, unsound way to translate a many-sorted FOL problem to unsorted FOL is to erase all the sorts and otherwise leave the problem unchanged. There are two main sound alternatives that encode the sort information. Sort *tags* are functions  $t_\sigma(X)$  that directly associate a term  $X$  with its sort  $\sigma$ . Sort *guards* are predicates  $g_\sigma(X)$  that check whether  $X$  has sort  $\sigma$  in the original problem. The formalized versions of these encodings follow the four steps sketched in Section 2.

**Full Erasure.** Full sort erasure is unsound but complete. What makes it interesting is that it is sound for the class of monotonic problems. By way of composition, it lies at the heart of the tag- and guard-based encodings. The theory prefix  $U$  distinguishes unsorted entities from their many-sorted counterparts.

**SIG:** The signature of the target unsorted problem has the same function and predicate symbols as the original signature but collapses the sorts into a single, implicit sort.

**TRANS:** The translation function  $e$  is the identity except that it forgets the sorts.

**SOUND:** For the soundness proof, a model of a monotonic problem is extended into a model that interprets all sorts infinitely, which in turn is transformed into an isomorphic “full” model that interprets all the sorts uniformly as  $\lambda d. \text{True}$  (i.e.,  $\forall \sigma. \forall d. \text{int}_S \sigma d$ ), from which it is easy to build an unsorted model for the  $e$ -translated problem:

$$\text{Mono\_Model} < \text{Infinite\_Model} < \text{Full\_Model} < \text{U\_Model}$$

The last step corresponds to Theorem 1 in Bouillaguet et al. [8] and, more approximately, to Lemma 1 in Claessen et al. [9]. Incidentally, the formalization revealed a flaw in Claessen et al.: Their main result holds, but not their Lemma 3.<sup>1</sup>

**COMPLETE:** The locale `Problem_UModel` combines a many-sorted problem and an unsorted model with domain  $D$  of the problem’s  $e$  translation. The unsorted model can be regarded as a many-sorted model in which every sort is interpreted as  $D$ .

**Protector-Based Encodings.** Claessen et al. observe that protectors, whether tags or guards, are not needed for terms with monotonic sorts. The sequel [3] advocates protecting only those variables that cause the monotonicity check to fail, to reduce clutter. Thus, for both tags and guards, three schemes are available: the traditional encoding, the lightweight version due to Claessen et al., and the “featherweight” version from the sequel. These are called  $\tilde{t}$ ,  $\tilde{t}?$ , and  $\tilde{t}??$  for tags and  $\tilde{g}$ ,  $\tilde{g}?$ , and  $\tilde{g}??$  for guards.

Consider the following fragment of a many-sorted problem, where  $S$  has sort  $\text{st}$ :

$$S \approx \text{on} \vee S \approx \text{off} \qquad \text{flip}(S) \not\approx S$$

<sup>1</sup> The flawed lemma states that whenever there exists a model  $\mathcal{M}$  where a monotonic sort  $\sigma$  is interpreted with a given cardinality, there exists for any larger cardinality  $k$  a model where  $\sigma$  has cardinality  $k$  and the other sorts have the same cardinalities as in  $\mathcal{M}$ . This proposition is invalid for  $k > \aleph_0$  because FOL problems can encode the constraint that there exists a bijection between two infinite, and hence monotonic, sorts  $\sigma$  and  $\tau$ , making it impossible to increase  $\sigma$ ’s cardinality without also increasing  $\tau$ ’s. This issue is independent of which of the two definitions of monotonicity is used. We discovered it at an early stage of the formalization as we were looking for a correct formulation of Löwenheim–Skolem for many-sorted FOL.



The traditional  $\tilde{t}$  encoding inserts tags around every subterm:

$$t_{st}(S) \approx t_{st}(\text{on}) \vee t_{st}(S) \approx t_{st}(\text{off}) \quad t_{st}(\text{flip}(t_{st}(S))) \not\approx t_{st}(S)$$

Since the sort  $st$  is not monotonic (its only models have cardinality 2), the  $\tilde{t}?$  encoding coincides with  $\tilde{t}$ . In contrast, the featherweight  $\tilde{t}??$  encoding tags only naked variables:

$$t_{st}(S) \approx \text{on} \vee t_{st}(S) \approx \text{off} \quad \text{flip}(S) \not\approx S$$

The  $\tilde{t}??$ -encoded problem is complemented by typing axioms that repair mismatches between tagged and untagged occurrences of well-sorted terms:

$$t_{st}(\text{on}) \approx \text{on} \quad t_{st}(\text{off}) \approx \text{off} \quad t_{st}(\text{flip}(S)) \approx \text{flip}(S)$$

For guards, the traditional and lightweight encodings  $\tilde{g}$  and  $\tilde{g}?$  protect each variable:

$$\neg g_{st}(S) \vee S \approx \text{on} \vee S \approx \text{off} \quad \neg g_{st}(S) \vee \text{flip}(S) \not\approx S$$

The featherweight encoding  $\tilde{g}??$  guards only naked variables:

$$\neg g_{st}(S) \vee S \approx \text{on} \vee S \approx \text{off} \quad \text{flip}(S) \not\approx S$$

The guard encodings are completed by the axioms  $g_{st}(\text{on})$ ,  $g_{st}(\text{off})$ , and  $g_{st}(\text{flip}(S))$ .

**General Encoding Procedure.** The full sort erasure encoding  $e$  is part of a two-stage procedure to encode any many-sorted FOL problem into unsorted FOL. The first stage makes the problem monotonic by introducing protectors (tags or guards). This corresponds to a sound and complete encoding of many-sorted FOL into itself; the soundness proofs rely on the monotonicity criteria. The second stage merges all the sorts using  $e$ , which is sound and complete for monotonic problems.

Tags and guards are formalized separately, but for a protector kind, the traditional, lightweight, and featherweight encodings are treated as instances of a single generalized encoding. Both generalized encodings are parameterized by a partition of sorts by level of protection, via disjoint predicates  $\text{prot}$ ,  $\text{protFw}$ ,  $\text{unprot} : 's \rightarrow \text{bool}$  indicating whether terms of a sort should be fully protected, protected in a featherweight fashion, or left unprotected. The last option is available only for sorts inferred monotonic by  $\triangleright$ .

**Tags.** The tag encoding builds on a datatype of extended function symbols containing the old symbols as well as a tag for each sort:

$$\text{datatype } (f, 's) \text{ efsym} = \text{Old } f \mid \text{Tag } 's$$

**SIG:** Signatures over the extended symbols treat the old function symbols as before. The new symbols  $\text{Tag } \sigma$  are unary operations of sort arity  $[\sigma]$  and result sort  $\sigma$ .

**TRANS:** The encoding function is specified as follows:

$$\begin{aligned} t(\text{Var } x) &= \begin{cases} \text{Var } x & \text{if } \text{unprot}(\text{sort } x) \\ \text{Fn } (\text{Tag } (\text{sort } x)) [\text{Var } x] & \text{otherwise} \end{cases} \\ t(\text{Fn } f \text{ } ts) &= t'(\text{Fn } f \text{ } ts) \\ t(\text{Pos } (\text{Eq } t_1 \text{ } t_2)) &= \text{Pos } (\text{Eq } (t \text{ } t_1) (t \text{ } t_2)) \\ t(\text{Neg } (\text{Eq } t_1 \text{ } t_2)) &= \text{Neg } (\text{Eq } (t' \text{ } t_1) (t' \text{ } t_2)) \\ t(\text{Pos } (\text{Pr } p \text{ } ts)) &= \text{Pos } (\text{Pr } p \text{ } (\text{map } t' \text{ } ts)) \\ t(\text{Neg } (\text{Pr } p \text{ } ts)) &= \text{Neg } (\text{Pr } p \text{ } (\text{map } t' \text{ } ts)) \end{aligned}$$

$$\begin{aligned}
t'(\text{Var } x) &= \begin{cases} \text{Fn}(\text{Tag}(\text{sort } x))[\text{Var } x] & \text{if } \text{prot}(\text{sort } x) \\ \text{Var } x & \text{otherwise} \end{cases} \\
t'(\text{Fn } f \text{ } ts) &= \begin{cases} \text{Fn}(\text{Tag}(\text{res } f))[\text{Fn}(\text{Old } f)(\text{map } t' \text{ } ts)] & \text{if } \text{prot}(\text{res } f) \\ \text{Fn}(\text{Old } f)(\text{map } t' \text{ } ts) & \text{otherwise} \end{cases}
\end{aligned}$$

The  $t$  function tags naked variables unless they are of an unprotected sort. The auxiliary function  $t'$  adds tags only for fully protected sorts; it is invoked on all subterms except naked variables.

The tag axioms  $\mathcal{A}_{x\Phi}$ —needed to repair mismatches between tagged and untagged terms in the featherweight encoding  $\tilde{t}??$ —have the form  $\text{Pos}(\text{Eq}(\text{Fn}(\text{Tag}(\text{res } f)[t]))t)$ , where  $t = \text{Fn}(\text{Old } f)(\text{map } \text{Var } xs)$  and  $xs$  is a list of distinct variables of sorts  $\text{arity}_F f$ , for all function symbols  $f$  such that  $\text{protFw}(\text{res } f)$ . The encoding of a problem is  $t\Phi = \{\text{map } t \text{ } c \mid c \in \Phi\} \cup \mathcal{A}_{x\Phi}$ .

**SOUND:** Given a model of the fixed problem  $\Phi$ , a model of  $t\Phi$  is obtained by extending it with interpreting tags as the identity functions.

**COMPLETE:** Completeness is more difficult. To convey a sense of the complexity, let us quote the informal proof, in which  $x$  stands for  $\tilde{t}?$  or  $\tilde{t}??$  ( $\tilde{t}$  is analogous to  $\tilde{t}?$ ) and  $\llbracket \Phi \rrbracket_x$  denotes the  $x$ -encoding of the NNF problem  $\Phi$  [4, §4.4]:

A model of  $\llbracket \Phi \rrbracket_x$  is *canonical* if all tag functions  $t_\sigma$  are interpreted as the identity. From a canonical model, we obtain a model of  $\Phi$  by leaving out  $t_\sigma$ . It then suffices to prove that whenever there exists a model  $\mathcal{M}$  of  $\llbracket \Phi \rrbracket_x$ , there exists a canonical model  $\mathcal{M}'$ .

For  $\tilde{t}?$ , values of a tagged type  $\sigma$  are systematically accessed through  $t_\sigma$ . Hence, we can safely permute the entries of the function table of each  $t_\sigma$  so that it is the identity for the values in its range. We then construct  $\mathcal{M}'$  by removing the domain elements for which  $t_\sigma$  is not the identity. It is a model by Lemma 4.13 [which states that substructures of NNF models are models if they preserve existential witnesses].

For  $\tilde{t}??$ , the construction must take possibly nonmonotonic types into account. No permutation is necessary for these thanks to the typing axioms, which ensure that the tag functions are the identity for well-typed terms. For each  $\sigma \not\vdash \Phi$ , we remove the model elements for which  $t_\sigma$  is not the identity. The typing axioms ensure that the substructure is well-defined: each tag function is the identity for at least one element and also for each element within the range of a non-tag function. The equations  $t_\sigma(X) \approx X$  generated for existential variables ensure that witnesses are preserved, as required by Lemma 4.13.

Relying on permutations is intuitive on paper, but in the proof assistant it is simpler to combine the permutation and the reduction to a canonical model:

$$\text{int}_F f \text{ } as = \begin{cases} \text{eint}_F(\text{Tag}(\text{res } f))[\text{eint}_F(\text{Old } f)(\text{map}_2 q(\text{arity}_F f) \text{ } as)] & \text{if } \text{prot}(\text{res } f) \\ \text{eint}_F(\text{Old } f)(\text{map}_2 q(\text{arity}_F f) \text{ } as) & \text{otherwise} \end{cases}$$

Here,  $\text{eint}_F$  denotes the  $\text{int}_F$  component of the fixed model of  $t\Phi$ , and  $\text{map}_2$  applies a binary function elementwise on parallel lists. The auxiliary function  $q$  maps a sort  $\sigma$

and an element  $d$  to  $d$  if either  $\text{unprot } \sigma$  or  $d$  is in the range of  $\text{eint}_F (\text{Tag } \sigma)$ ; otherwise, it maps  $\sigma, d$  to  $\text{eint}_F (\text{Tag } \sigma) d$ . The proof that the resulting structure is a model of the original problem  $\Phi$  involves defining suitable back-and-forth functions between the two structures. Finally, proving monotonicity of  $\text{t } \Phi$  is reduced to showing that the first criterion always succeeds on the translated problem:  $\text{Problem}_\Phi < \text{Problem\_Crit1}_{\text{t}\Phi}$ .

**Guards.** The guard encoding requires extending the signature with guard predicates:

`datatype ('p, 's) epsym = Old 'p | Guard 's`

Each symbol  $\text{Guard } \sigma$  has arity  $[\sigma]$  and contributes axioms to the translated problem.

The soundness proof extends models of  $\Phi$  into models of  $\text{g } \Phi$  by interpreting the guard predicates as true everywhere. The completeness part is easier for guards than for tags. A canonical model is one where all guard predicates are interpreted as true everywhere. The proof handles the three levels of protection uniformly, reflecting the more uniform nature of  $\tilde{\text{g}}??$ —there are no counterparts to the “typing axioms that repair mismatches between tagged and untagged occurrences of well-sorted terms” of  $\tilde{\text{t}}??$ .

Monotonicity is proved using the second criterion, with the extension policy  $\text{C}$  for the predicates  $\text{Old } p$  and  $\text{F}$  for the distinguished symbols  $\text{Guard } \sigma$ . This is a departure from the informal proof, which inlines the model extension argument without appealing to the monotonicity criterion.

## 6 First-Order Logic with Quantifiers

This and the next two sections are concerned with lifting the results presented in the previous sections to negation normal form and structures with arbitrarily large domains.

The locales for quantified FOL formulas in NNF are either the same or similar to those for CNF; the theory prefix  $\text{Q}$  is used for disambiguation (e.g.,  $\text{Q.Model}$ ). No cardinality assumption is made about the universe. Terms and atoms are as for CNF, but formulas can nest positive connectives and quantifiers arbitrarily.

The following declaration gives an approximation of the syntactic category of formulas. The actual type identifies them modulo  $\alpha$ -equivalence (variable renaming):

`datatype ('s, 'f, 'p) fm =  
 Pos (('f, 'p) atm) | Conj (('s, 'f, 'p) fm) (('s, 'f, 'p) fm) | All 's var (('s, 'f, 'p) fm) |  
 Neg (('f, 'p) atm) | Disj (('s, 'f, 'p) fm) (('s, 'f, 'p) fm) | Ex 's var (('s, 'f, 'p) fm)`

The proper formal management of binding syntax modulo  $\alpha$ -equivalence is a topic of extensive research in  $\lambda$ -calculus and programming languages. FOL syntax poses similar challenges. In particular, substitution and its interplay with the semantics is difficult to handle rigorously; for example, a standard textbook [15] dedicates dozens of lemmas to these preliminaries, with rough proof sketches. Many of these refer to properties of any syntax with static bindings, falling under the scope of a general metatheory of syntax formalized by Popescu et al. [22–24]. A prominent feature of this framework—distinguishing it from the more established Nominal Isabelle [13], based on nominal logic [20]—is that it is centered around the notion of substitution:

- The framework defines substitution, including parallel and unary variants, and provides a large collection of basic facts about the interaction of substitution with free variables and the other operators.

- It provides a recursor for defining operators that are directly compositional with substitution. (In contrast, the nominal logic recursor targets compositionality with permutations, a less useful concept.)

This unconventional focus is appropriate: Substitution is without doubt the central syntactic operator in logics and type systems.

Another main feature is the facilitation of semantic interpretation of syntax, which is problematic in frameworks optimized for manipulating finitary syntax. For example, Pitts encounters “a really nontrivial freshness condition on binders” [21, §6.3] he needs to discharge in the context of applying the nominal recursor to interpret the  $\lambda$ -calculus in a semantic domain. This feature is illustrated below for interpreting FOL syntax.

The framework requires the user to provide semantic domains—for FOL, types  $\mathcal{T}$ ,  $\mathcal{A}$ , and  $\mathcal{F}$  for interpreting terms, atoms, and formulas—as well as first-order operations corresponding to the non-binding constructors other than for variables (e.g.,  $\text{FN} : 'f \rightarrow \mathcal{T} \text{ list} \rightarrow \mathcal{T}$ ) and second-order operations corresponding to the binders:  $\text{ALL} : 's \rightarrow (\mathcal{T} \rightarrow \mathcal{F}) \rightarrow \mathcal{F}$  and  $\text{EX} : 's \rightarrow (\mathcal{T} \rightarrow \mathcal{F}) \rightarrow \mathcal{F}$ .

In exchange, the framework produces the functions  $\text{int}_{\text{Tm}} : \text{tm} \rightarrow (\text{var} \rightarrow \mathcal{T}) \rightarrow \mathcal{T}$ ,  $\text{int}_{\text{At}} : \text{atm} \rightarrow (\text{var} \rightarrow \mathcal{T}) \rightarrow \mathcal{A}$ , and  $\text{int}_{\text{Fm}} : \text{fm} \rightarrow (\text{var} \rightarrow \mathcal{T}) \rightarrow \mathcal{F}$  that interpret syntax in the semantic domains. They map variables according to a valuation  $\xi$ . They map the action of non-binding constructors to that of the corresponding semantic operators, and similarly for binding constructors but in a valuation-sensitive way. For example:

$$\begin{aligned} \text{int}_{\text{Tm}} (\text{Var } x) \xi &= \xi x \\ \text{int}_{\text{Tm}} (\text{Fn } f \text{ ts}) \xi &= \text{FN } f (\text{map } (\lambda t. \text{int}_{\text{Tm}} t \xi) \text{ ts}) \\ \text{int}_{\text{At}} (\text{Eq } t_1 t_2) \xi &= \text{EQ } (\text{int}_{\text{Tm}} t_1 \xi) (\text{int}_{\text{Tm}} t_2 \xi) \\ \text{int}_{\text{Fm}} (\text{Conj } \varphi_1 \varphi_2) \xi &= \text{CONJ } (\text{int}_{\text{Fm}} \varphi_1 \xi) (\text{int}_{\text{Fm}} \varphi_2 \xi) \\ \text{int}_{\text{Fm}} (\text{All } \sigma x \varphi) \xi &= \text{ALL } \sigma (\lambda d. \text{int}_{\text{Fm}} \varphi \xi[x \mapsto d]) \end{aligned}$$

where  $\xi[x \mapsto d]$  denotes the function that maps  $x$  to  $d$  and otherwise coincides with  $\xi$ . So far, this looks like the standard interpretation of binding syntax in a semantic domain, except that here the recursive definition is modulo  $\alpha$ -equivalence (which is a priori difficult to achieve in a proof assistant). The framework also derives compositionality of substitution w.r.t. valuation update and obliviousness of the interpretation w.r.t. fresh variables in a systematic, FOL-agnostic way:

$$\begin{aligned} \text{int}_{\text{Fm}} \varphi[t/x] \xi &= \text{int}_{\text{Fm}} \varphi \xi[x \mapsto \text{int}_{\text{Tm}} t \xi] \\ \text{int}_{\text{Fm}} \varphi \xi &= \text{int}_{\text{Fm}} \varphi \xi' \quad \text{if } \xi \text{ and } \xi' \text{ differ only on variables fresh for } \varphi \end{aligned}$$

In the first equation,  $\varphi[t/x]$  denotes capture-free substitution of  $t$  for  $x$  in  $\varphi$ .

A many-sorted structure  $(\text{int}_{\text{S}}, \text{int}_{\text{F}}, \text{int}_{\text{P}})$  can be organized as a semantic domain by taking  $\mathcal{T} = \omega$ ,  $\mathcal{A} = \mathcal{F} = \text{bool}$ ,  $\text{FN} = \text{int}_{\text{F}}$ ,  $\text{EQ} = (=)$ ,  $\text{CONJ} = (\wedge)$ ,  $\text{ALL } \sigma P = (\forall d. \text{int}_{\text{S}} \sigma d \longrightarrow P d)$ , and so on. This yields the recursive equations

$$\begin{aligned} \llbracket \text{Var } x \rrbracket_{\xi} &= \xi x \\ \llbracket \text{Fn } f \text{ ts} \rrbracket_{\xi} &= \text{int}_{\text{F}} f (\text{map } (\lambda t. \llbracket t \rrbracket_{\xi}) \text{ ts}) \\ \models_{\xi} \text{Eq } t_1 t_2 &\longleftrightarrow \llbracket t_1 \rrbracket_{\xi} = \llbracket t_2 \rrbracket_{\xi} \\ \models_{\xi} \text{Conj } \varphi_1 \varphi_2 &\longleftrightarrow \models_{\xi} \varphi_1 \wedge \models_{\xi} \varphi_2 \\ \models_{\xi} \text{All } \sigma x \varphi &\longleftrightarrow \forall d. \text{int}_{\text{S}} \sigma d \longrightarrow \models_{\xi[x \mapsto d]} \varphi \end{aligned}$$

which characterize term interpretation (with  $\llbracket t \rrbracket_\xi = \text{int}_{\text{tm}} t \xi$ ), atom satisfaction ( $\models_\xi a = \text{int}_{\text{at}} a \xi$ ), and formula satisfaction ( $\models_\xi \varphi = \text{int}_{\text{fm}} \varphi \xi$ ). These functions are defined in the `Q.Structure` locale. The framework also produces the substitution lemma  $\models_\xi \varphi[t/x] \longleftrightarrow \models_{\xi[x \mapsto \llbracket t \rrbracket_\xi]} \varphi$ . In the next section, the notations  $\models \varphi$  and  $\models \Phi$  abbreviate  $\forall \xi. \models_\xi \varphi$  and  $\forall \varphi \in \Phi. \models \varphi$ . The structure can also be made explicit—e.g.,  $(\text{int}_s, \text{int}_f, \text{int}_p) \models_\xi \varphi$ .

If the orientation toward substitution is the main strength of the framework, its main weakness is the lack of automation. For each desired binding syntax type, users must currently instantiate the general theorems manually, much like mathematicians do routinely when applying universal algebra to groups or rings. The instantiation is tedious due to the large number of theorems. Despite the availability of “template files,” this process can take days and thousands of lines of proof text. Automation in the form of a definitional package, which would provide the basic convenience expected by users of Nominal Isabelle (while supporting substitution natively), remains for future work.

## 7 Classical Metatheorems

The lifting argument from countable CNF structures to unbounded NNF structures (Section 8) relies on clausification and Löwenheim–Skolem for many-sorted FOL with equality. Earlier formalizations focus on unsorted FOL without equality [2, 11, 25]. Sorts and equality are tedious to formalize, and they often fail to reward the formalizer with deep logical insight, but they are central to monotonicity and sort encodings.

**Clausification.** The translation of a finite quantified problem into clausal form involves skolemizing all the existentially quantified variables into function symbols that take the universally quantified variables in scope as arguments. Skolemization is surprisingly difficult to treat formally; for example, Harrison [11] claims that it poses greater challenges than completeness. On the positive side, clausification can be seen as an instance of the general semantic interpretation principle introduced in Section 6.

The definition of clausification and its soundness and completeness proof follow the four-step institutional approach.

**SIG:** Skolemization introduces new function symbols  $\text{Sko } \sigma s \ x$ , built from a list of sorts  $\sigma s$  (specifying the arity) and a variable name  $x$ , while preserving the sorts of  $\Sigma$ -symbols:

$$\text{datatype } 'f \text{efsym} = \text{Old } 'f \mid \text{Sko } ('s \text{ list}) \text{ var}$$

**TRANS:** The clausification function  $\text{cls}$  takes a  $\Sigma$ -formula  $\varphi$ , an environment  $\rho : \text{var} \rightarrow \text{tm}$ , a list of universal variables  $vs$ , and a set of fresh variables  $V$  as arguments. In addition to massaging the connectives, it replaces existential variables by new symbols that depend on  $vs$ , replaces bound universal variables by fresh variables from  $V$ , and substitutes free variables according to  $\rho$  to produce a  $\Sigma'$ -clause.

The characteristic equations for  $\text{cls}$  are obtained by instantiating the semantic interpretation principle with  $\mathcal{T} = \text{tm}$ ,  $\mathcal{A} = \text{atm}$ , and  $\mathcal{F} = \text{var list} \rightarrow \text{var set} \rightarrow \text{fm}$ , taking suitable operators on these domains, and letting  $\text{cls}$  be  $\text{int}_{\text{fm}}$ . The interesting cases are

$$\begin{aligned} \text{cls } (\text{All } \sigma \ x \ \varphi) \ \rho \ vs \ V &= \text{cls } \varphi \ \rho[x \mapsto \text{Var } v] \ (v \# vs) \ (V \setminus \{v\}) \\ \text{cls } (\text{Ex } \sigma \ x \ \varphi) \ \rho \ vs \ V &= \text{cls } \varphi \ \rho[x \mapsto \text{Fn } f \ (\text{map Var } vs)] \ vs \ (V \setminus \{v\}) \end{aligned}$$

where  $v \in V$  is some variable of sort  $\sigma$  and  $f = \text{Sko}(\text{map sort } vs) v$  is the Skolem function symbol, which is applied to the universal variables  $vs$ . For closed formulas, clausification is defined as  $\text{clausify } \varphi = \text{cls } \varphi \rho [] \text{ UNIV}$  for some irrelevant choice of  $\rho$ .

As a simple example, let  $\varphi = \text{All } \sigma x (\text{Ex } \tau y (\text{Eq}(\text{Var } x) (\text{Var } y)))$ , let  $v_1, v_2$  be the variables picked from  $\text{UNIV}$  and  $\text{UNIV} \setminus \{v_1\}$ , and let  $f = \text{Sko}[\sigma] v_2$ . Then

$$\begin{aligned} & \text{clausify } \varphi \\ &= \text{cls } \varphi \rho [] \text{ UNIV} \\ &= \text{cls } (\text{Ex } \tau y (\text{Eq}(\text{Var } x) (\text{Var } y))) \rho[x \mapsto \text{Var } v_1] [v_1] (\text{UNIV} \setminus \{v_1\}) \\ &= \text{cls } (\text{Eq}(\text{Var } x) (\text{Var } y)) \rho[x \mapsto \text{Var } v_1, y \mapsto \text{Fn } f [\text{Var } v_1]] [v_1] (\text{UNIV} \setminus \{v_1, v_2\}) \\ &= \text{Eq}(\text{Var } v_1) (\text{Fn } f [\text{Var } v_1]) \end{aligned}$$

**SOUND:** Soundness is proved in the Structure locale, which fixes a  $\Sigma$ -structure  $(\text{int}_S, \text{int}_F, \text{int}_P)$ . The “Skolem model” predicate  $\text{skmod } \varphi \rho vs V \text{ eint}_F \text{ eint}'_F$  transforms, for each valuation  $\xi : \text{var} \rightarrow 'u$ , an extended structure  $\text{eint}_F$  such that  $\models_\xi \text{cls } \varphi \rho vs V$  into an extended structure  $\text{eint}'_F$  such that  $\models_{\xi \diamond \rho} \varphi$ , where  $\diamond$  composes valuations with environments. The introduction rules of  $\text{skmod}$  emulate  $\text{cls}$ ’s equations; for example,

$$\frac{\text{skmod } \varphi \rho[x \mapsto \text{Fn } f (\text{map Var } vs)] vs (V \setminus \{v\}) \text{ eint}_F[f \mapsto F] \text{ eint}'_F}{\text{skmod } (\text{Ex } \sigma x \varphi) \rho vs V \text{ eint}_F \text{ eint}'_F}$$

where  $v \in V$  and  $F : 'u \text{ list} \rightarrow 'u$  is a suitable interpretation for the Skolem symbol  $f$ , defined so that  $F us$  gives an arbitrary  $u$  such that  $(\text{int}_S, \text{int}_F, \text{int}_P) \models_{\xi \diamond \rho[x \mapsto u]} \varphi$ , where  $\xi$  maps  $vs$  to  $us$  elementwise. The  $\text{skmod}$  relation is total on the last argument. For closed formulas  $\varphi$  such that  $(\text{int}_S, \text{int}_F, \text{int}_P) \models \varphi$ , starting with an extension  $\text{eint}_F$  of  $\text{int}_F$ ,  $\text{skmod}$  yields  $\text{eint}'_F$  such that  $(\text{int}_S, \text{eint}'_F, \text{int}_P) \models \text{clausify } \varphi$ . Thus, if  $\varphi$  has a model, then  $\text{clausify } \varphi$  also has a model.

For problems, we define  $\text{clausify } \Phi = \text{clausify } (\bigwedge \Phi)$ , where  $\bigwedge \Phi$  is the conjunction of all formulas in  $\Phi$ , which must be finite. The locale  $\text{Q.Model}$  fixes  $\Phi$  and a model, which is also a model of the formula  $\bigwedge \Phi$ . By soundness of  $\text{clausify}$  on closed formulas, this yields a model of  $\text{clausify } \Phi$ .

**COMPLETE:** For completeness, it suffices to show that the backward structure translation of a model of  $\text{clausify } \Phi$  is a model of  $\Phi$ . This is straightforward.

**Löwenheim–Skolem.** The proof of the downward Löwenheim–Skolem theorem is based on a formalization of a complete inference system, described in a separate paper [7]. In the  $\text{Q.Model}$  locale, which fixes a problem and model, it constructs a syntactic Henkin model. Since this model has a countable universe, there exists an isomorphic copy on  $\omega$  (the countably infinite universe fixed throughout Sections 4 and 5). This yields  $\text{Q.Model}_{\iota_u} < \text{Q.Model}_\omega$ .

Using the obvious sound and complete embedding  $\text{embed}$  of CNF problems into NNF problems, it is possible to transfer the Löwenheim–Skolem theorem to CNF:

$$\text{Model}_{\iota_u, \Phi} < \text{Q.Model}_{\iota_u, \text{embed } \Phi} < \text{Q.Model}_{\omega, \text{embed } \Phi} < \text{Model}_{\omega, \Phi}$$

To summarize the results of this section:

**Theorem 2.** *An NNF problem  $\Phi$  has a model iff  $\text{clausify } \Phi$  has a model.*

**Theorem 3.** *An NNF problem has a model iff it has a countable model.*

## 8 Lifting to Arbitrary Structures and Formulas with Binders

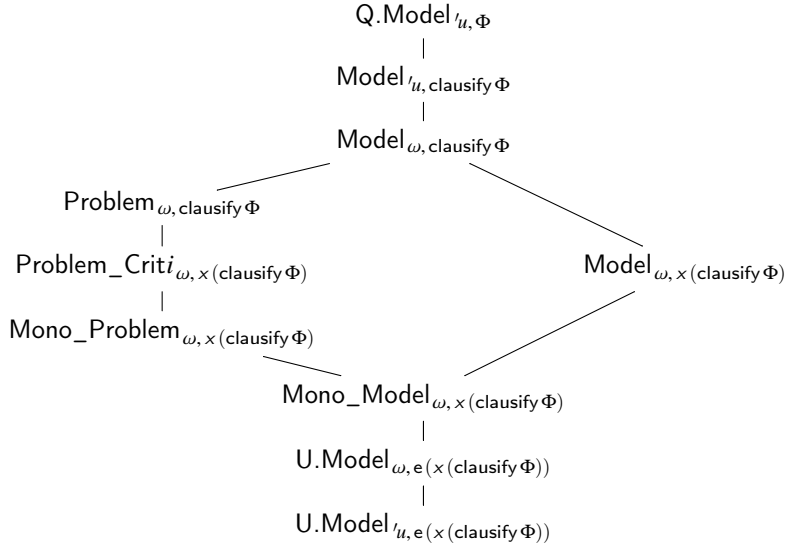
The focus on clausal form and countable structures is a useful simplification, but it is not faithful to the NNF-based paper proof [3] (or to the implementation in Sledgehammer). Thanks to a lifting argument that relies on clausification and Löwenheim–Skolem, the final results are free of such restrictions.

Figure 1 shows how the results are connected. Starting at the top with a satisfiable quantified problem  $\Phi$ , the problem is first clausified, then by Löwenheim–Skolem it is countably satisfiable (by taking  $u = \omega$ ). On the left-hand side, the clausified problem is further encoded using tags or guards ( $x \in \{t, g\}$ ) and shown to pass one of the monotonicity criteria ( $i = 1$  for  $t$  and  $2$  for  $g$ ), meaning it is monotonic. On the right-hand side, the encoded problem is satisfiable. Merging the two branches yields a monotonic satisfiable problem, whose erasure is a satisfiable unsorted problem. Since every translation step is also shown complete, the right-hand side can also be traversed bottom-up, producing a model of the original problem from a model of the translated unsorted problem. The overall translation is thus sound and complete.

**Theorem 4.** *Given  $x \in \{t, g\}$  and a finite many-sorted NNF problem  $\Phi$ , let  $\Phi'$  be the unsorted CNF problem  $e(x(\text{clausify } \Phi))$ , i.e., the sort-erased  $x$ -translated clausified  $\Phi$ .*

- (1) *For each model  $\mathcal{M}$  of  $\Phi$  (forming together with  $\Phi$  an instance of  $\text{Q.Model}_\Phi$ ), there exists a model  $\mathcal{M}'$  of  $\Phi'$  (forming together with  $\Phi'$  an instance of  $\text{U.Model}_{\Phi'}$ ).*
- (2) *Conversely, for every model of  $\Phi'$ , there exists a model of  $\Phi$ .*

The formal proof puts together many constructions and results of independent interest, notably soundness of the monotonicity criteria (Theorem 1), soundness and completeness of clausification (Theorem 2), and downward Löwenheim–Skolem (Theorem 3).



**Figure 1.** The verified translation pipeline

## 9 Conclusion

This paper describes a framework and a methodology for formalizing applications of many-sorted first-order logic while acting as a companion to recent papers on sort encodings [3, 8, 9]. To readers from the proof assistant community, it also provides a contribution to the ongoing binder representation debate. And to readers rooted in algebraic methods, it shows a practical application of the theory of institutions in a context where the translation functions cannot be assumed to be uniform.

The formalization widely reaffirmed already proved results. On one occasion, it revealed a flaw in a published lemma (Lemma 3 of Claessen et al. [9]). It also helped detect mistakes in a subsequent paper proof [4] before it reached any readers. The work provided the opportunity to rethink the proof; for example, the generalized monotonicity concept, in terms of sets of sorts, arose during the formalization.

A potential practical benefit of this work is connected to step-by-step proof reconstruction. Although the encodings are sound, the inferences in a machine-generated proof may violate the sort discipline, resulting in failures in Sledgehammer's proof replay. In future work, we want to investigate the feasibility of connecting the soundness proofs of the encodings with a verified checker for unsorted FOL proofs.

The advantages of machine-checked metatheory are well known from programming language research, where papers are often accompanied by formal developments and proof assistants have made it into the classroom. Paradoxically, in the automated reasoning community, we have not been very enthusiastic about formalizing our own results. This paper reported on some steps we have taken to address this.

**Acknowledgement.** We thank Tobias Nipkow for making this work possible. Jesper Bengtson, Nicholas Smallbone, Mark Summerfield, Dmitriy Traytel, and several anonymous reviewers suggested improvements to earlier versions of this paper. The research was supported by the Deutsche Forschungsgemeinschaft (DFG) projects Security Type Systems and Deduction (grant Ni 491/13-1), part of the program Reliably Secure Software Systems (RS<sup>3</sup>, Priority Program 1496), and Hardening the Hammer (grant Ni 491/14-1). The authors are listed in alphabetical order.

## References

- [1] Ballarín, C.: Locales: A module system for mathematical theories. *J. Autom. Reasoning*, to appear
- [2] Berghofer, S.: First-order logic according to Fitting. In: Klein, G., Nipkow, T., Paulson, L. (eds.) *Archive of Formal Proofs*. <http://afp.sf.net/entries/FOL-Fitting.shtml> (2007)
- [3] Blanchette, J.C., Böhme, S., Popescu, A., Smallbone, N.: Encoding monomorphic and polymorphic types. In: Piterman, N., Smolka, S. (eds.) *TACAS 2013. LNCS*, vol. 7795, pp. 493–507. Springer (2013)
- [4] Blanchette, J.C., Böhme, S., Popescu, A., Smallbone, N.: Encoding monomorphic and polymorphic types. Tech. report associated with TACAS 2013 paper [3], [http://www21.in.tum.de/~blanchet/enc\\_types\\_report.pdf](http://www21.in.tum.de/~blanchet/enc_types_report.pdf) (2013)
- [5] Blanchette, J.C., Krauss, A.: Monotonicity inference for higher-order formulas. *J. Autom. Reasoning* 47(4), 369–398 (2011)



- [6] Blanchette, J.C., Popescu, A.: Formal development associated with this paper. [http://www21.in.tum.de/~popescua/foL\\_devel.zip](http://www21.in.tum.de/~popescua/foL_devel.zip) (2013)
- [7] Blanchette, J.C., Popescu, A., Traytel, D.: Coinductive pearl: Modular first-order logic completeness, submitted, <http://www21.in.tum.de/~blanchet/compl.pdf>
- [8] Bouillaguet, C., Kuncak, V., Wies, T., Zee, K., Rinard, M.C.: Using first-order theorem provers in the Jahob data structure verification system. In: Cook, B., Podelski, A. (eds.) VMCAI 2007. LNCS, vol. 4349, pp. 74–88. Springer (2007)
- [9] Claessen, K., Lillieström, A., Smallbone, N.: Sort it out with monotonicity—Translating between many-sorted and unsorted first-order logic. In: Bjørner, N., Sofronie-Stokkermans, V. (eds.) CADE-23. LNAI, vol. 6803, pp. 207–221. Springer (2011)
- [10] Goguen, J.A., Burstall, R.M.: Institutions: Abstract model theory for specification and programming. J. ACM 39(1), 95–146 (1992)
- [11] Harrison, J.: Formalizing basic first order model theory. In: Grundy, J., Newey, M.C. (eds.) TPHOLs '98. LNCS, vol. 1479, pp. 153–170. Springer (1998)
- [12] Harrison, J.: Towards self-verification of HOL Light. In: Furbach, U., Shankar, N. (eds.) IJCAR 2006. LNCS, vol. 4130, pp. 177–191. Springer (2006)
- [13] Huffman, B., Urban, C.: Proof pearl: A new foundation for Nominal Isabelle. In: Kaufmann, M., Paulson, L.C. (eds.) ITP 2010. LNCS, vol. 6172, pp. 35–50. Springer (2010)
- [14] Kammüller, F., Wenzel, M., Paulson, L.C.: Locales—A sectioning concept for Isabelle. In: Bertot, Y., Dowek, G., Hirschowitz, A., Paulin, C., Théry, L. (eds.) TPHOLs '99. LNCS, vol. 1690, pp. 149–166. Springer (1999)
- [15] Monk, J.D.: Mathematical Logic. Springer (1976)
- [16] Myreen, M.O., Davis, J.: A verified runtime for a verified theorem prover. In: van Eekelen, M.C.J.D., Geuvers, H., Schmaltz, J., Wiedijk, F. (eds.) ITP 2011. LNCS, vol. 6898, pp. 265–280. Springer (2011)
- [17] Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL: A Proof Assistant for Higher-Order Logic, LNCS, vol. 2283. Springer (2002)
- [18] Paulson, L.C., Blanchette, J.C.: Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In: Sutcliffe, G., Ternovska, E., Schulz, S. (eds.) IWIL-2010 (2010)
- [19] Pfenning, F., Elliott, C.: Higher-order abstract syntax. In: Wexelblat, R.L. (ed.) PLDI '88. pp. 199–208. ACM (1988)
- [20] Pitts, A.M.: Nominal logic, a first order theory of names and binding. Inf. Comput. 186(2), 165–193 (2003)
- [21] Pitts, A.M.: Alpha-structural recursion and induction. J. ACM 53(3), 459–506 (2006)
- [22] Popescu, A.: Contributions to the Theory of Syntax with Bindings and to Process Algebra. Ph.D. thesis, C.S. Dept., University of Illinois (2010)
- [23] Popescu, A., Gunter, E.L.: Recursion principles for syntax with bindings and substitution. In: Chakravarty, M.M.T., Hu, Z., Danvy, O. (eds.) ICFP 2011. pp. 346–358. ACM (2011)
- [24] Popescu, A., Gunter, E.L., Osborn, C.J.: Strong normalization of System F by HOAS on top of FOAS. In: LICS 2010. pp. 31–40. IEEE (2010)
- [25] Ridge, T., Margetson, J.: A mechanically verified, sound and complete theorem prover for first order logic. In: Hurd, J., Melham, T.F. (eds.) TPHOLs 2005. LNCS, vol. 3603, pp. 294–309. Springer (2005)
- [26] Shankar, N.: Metamathematics, Machines, and Gödel's Proof, Cambridge Tracts in Theoretical Computer Science, vol. 38. Cambridge University Press (1994)
- [27] Sutcliffe, G.: The 6th IJCAR automated theorem proving system competition—CASC-J6. AI Comm. 26(2), 211–223 (2013)
- [28] Tinelli, C., Zarba, C.: Combining decision procedures for sorted theories. In: Alferes, J., Leite, J. (eds.) JELIA 2004. LNCS, vol. 3229, pp. 641–653. Springer (2004)